

Checking properties after shallow parsing

Lúisa Coheur¹ and Nuno J. Mamede²

¹ L²F – INESC-ID Lisboa/IST/GRIL

² L²F – INESC-ID Lisboa/IST

Spoken Language Systems Lab

Rua Alves Redol 9, 1000-029 Lisboa, Portugal

{Luisa.Coheur,Nuno.Mamede}@inesc-id.pt

Resumo Tira-Teimas is a program that verifies if a shallow parsed text satisfies a set of properties from the 5P paradigm. Tira-Teimas is implemented in XSLT and it is automatically generated by another program called Tira-Teimas Translator, which is also implemented in XSLT.

1 Introduction

The 5P paradigm [1,2,3] uses a set of properties to describe natural language syntax and Hagège in [3] made a precise description of Portuguese nominal phrases using these 5P properties. Also, she developed a shallow parser prototype called AF [3,4], and enriched the structures used by AF with the information expressed within these properties. Nevertheless, AF's information structures are less expressive than 5P properties. Therefore, it is not sure that the sequences identified by AF verify the whole set of 5P properties.

Tira-Teimas was then built in order to check if there is any inconsistency between the information from 5P properties and AF's results. Therefore, Tira-Teimas is a program that verifies if a shallow parsed text satisfies or not a set of properties as it checks if each sequence identified by AF satisfies (or not) the 5P properties.

The following chapter briefly describes 5P properties, AF's structures, and emphasize the gap between these formalisms. Then, in Section 3 it is explained how Tira-Teimas was built. An auxiliary tool, called Tira-Teimas translator is introduced in Section 4. Section 5 shows some results and we conclude the document with the usual conclusions and future work.

2 5P properties, blocks and leaves

2.1 5P properties: a brief overview

In the following we briefly describe 5P properties (see [1,2] or [3] for extra details). This presentation is illustrated with examples from [3], where a detailed

description of the syntax of Portuguese nominal phrases (excluding coordination) can be found.

5P properties are divided into three sets: existence properties (vocabulary, uniqueness, nucleus, exigency and exclusion), linearity and arrow properties.

Vocabulary (V) This property allows to declare the elements (models or categories) occurring within a model. As an example

$$V_{Nn} = \{n, \text{adj}, \text{card_p}, \dots\}$$

indicates that inside a nuclear nominal phrase (Nn), we can have elements labeled by the elements from V_{Nn} or labeled by elements that represent a category subsumed by an element from V_{Nn} . In this particular example, the category labeled n(oun) subsumes the category for common nouns (singular, plural, ...) as well as proper nouns (of type 1, 2, 3 and 4). Therefore, elements from these categories can occur inside an Nn. The same happens with the categories subsumed by adj(ective). The category labelled card_p does not subsume any other category, meaning that an element labeled card_p can appear inside an Nn.

Uniqueness (Un) With this property we are able to identify the elements that cannot occur more than once in a certain model. As all example

$$Un_N = \{Nn\}$$

indicates that inside a nominal phrase (N), there is, at most, one nuclear nominal phrase Nn.

Nucleus (Nu) This property declares which are the elements that can be the nucleus of the phrase³. As so,

$$Nu_N = \{Nn\}$$

means that only a nominal nuclear phrase can be the nucleus of a nominal phrase.

Exigency (\Rightarrow) Exigency properties allow to declare that a certain element occurs in a model only if another element also occurs in it. As an example

$$\{\text{mesmo}\} \Rightarrow_{Nn} \{\text{dem}\} \mid \{\text{artd}\},$$

means that inside a nuclear nominal phrase, if there is an element labeled mesmo, then an element labeled dem(onstrative) or artd (definite article) must also occur

³ In each model there is one and only one nucleus.

(ex: *o mesmo, este mesmo carro, esse mesmo rapaz*)⁴.

We can have the following extra-notation for an exigency property:

$a \Rightarrow_X \text{NIL}$

meaning that **a** must be alone in a model labelled **X** (it excludes everything).

Exclusion (\nRightarrow) These properties permit to declare that an element excludes another element within a model. For example, with

$\{\text{ord}\} \nRightarrow_{Nn} \{\text{nc2}\}$

we have that an **ord**(inal) and a common noun of type 2 (**nc2**) exclude mutually inside a **Nn** (ex: *A terceira água*)⁵.

Linearity ($<$) These properties are responsible to declare linearity relations between the elements of a model. Consider

$\text{det} <_{Nn} \text{poss}$

It indicates that if a **det**(erminer) and **poss**(essive) occur in **Nn**, then the **det** precedes the **poss** (ex: *o meu livro*)⁶.

We say that a model verifies these properties if it does not disrespects any of them. As an example, *meu o livro*⁷ does not satisfies the presented properties because *meu* (**poss**) precedes the **det** *o*, desrespecting the linearity property.

2.2 Blocks and leaves

Hagège developed a surface analyser AF (the leaves algorithm), that receives, as input, two data structures called leaves and blocks. Briefly, each model has a block associated and each category has a leaf. Both leaves or blocks indicate:

- which is the model where they are in (the upper model);
- if they can always(1)/never(O)/sometimes(2) start the upper model;
- if they can always(1)/never(O)/sometimes(2) end the upper model;
- which are the elements that can appear after them in the upper model.

That is, each leaf and each block have the following structures, respectively:

⁴ Respectively, *the same, this same car, the same boy*.

⁵ *The third water*.

⁶ *My book*.

⁷ **my the book*.

leaf(-, category, upper model, start, end, next elements)

block(current model, upper model, start, end, next elements)

where

start, end $\in \{0, 1, 2\}$.

2.3 The gap

The problem is that the blocks and the leaves do not have the same expressiveness of 5P properties. As an example, consider the set of sequences labelled X and described by the following properties (suppose that a , b and c are categories from the vocabulary, and to simplify, that they are unique and that they can all be the nucleus of a model):

- (1) $a \Rightarrow_X b$
- (2) $a \not\Rightarrow_X c$
- (3) $a <_X b$

The following models (labeled X) respect these properties:

(c)_X
(b)_X
(a b)_X
(b c)_X
(c b)_X

However, the following ones do not respect them:

(a)_X
(b a)_X
(a c)_X
(a b c)_X
...

In order to accept the first set of sequences, AF needs the following leaves:

leaf(-, a, X, 1, 0, [b]);
leaf(-, b, X, 2, 2, [c]);
leaf(-, c, X, 2, 2, [b]).

because

- if a precedes and exiges b , then a always starts and never ends a model;

- as nothing is said about the linearity between **b** and **c**, **c** can precede **b** or appear after it, meaning that **b** may start or end (or not) the model;
- *mutatis mutandis* for **c**.

Therefore, with these leaves, the following sequence that was not verifying the 5P properties

$(a\ b\ c)_x$

would be consider an **X** model by AF, because:

- **a** may start a model;
- after **a** a **b** can occur;
- **c** can follow **a**;
- **c** can end a model.

3 Building Tira-Teimas

3.1 Choosing XSLT

Susana [5] is, in rough terms, a new implementation of AF. As Susana’s output is in XML [6], we decided to use XSLT (Extensible Stylesheet transformation) [7] in order to implement Tira-Teimas. XSLT is a tree-oriented transformation language, that uses stylesheets containing template rules and transforms XML into XML, HTML, txt, etc. As so, Tira-Teimas is nothing more than a stylesheet containing the 5P properties. It picks a shallow parsed text in XML, looks inside the detected sequences, and returns another shallow parsed text in XML, where the sequences not respecting 5P properties are identified and marked (see Fig. 1).

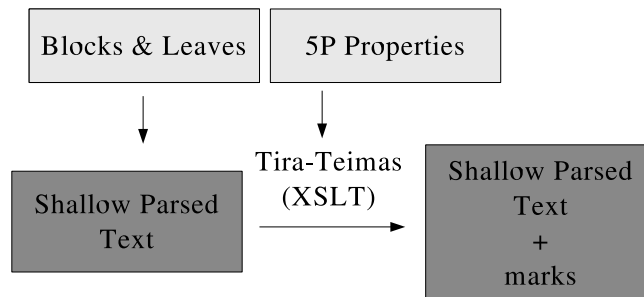


Figura 1. Tira-Teimas

3.2 Coding 5P properties in XSLT

In the following, we explain how 5P properties are coded in XSLT.

Exigency properties

Exigency properties have the following general syntax:

$$\text{Exig}_i: \{a_1, \dots, a_n\} \Rightarrow_M \{b_1, \dots, b_m\} \mid \dots \mid \{c_1, \dots, c_k\}$$

meaning that if the symbols a_1, \dots, a_n , occur in a model labelled M, then

$$\{b_1, \dots, b_m\} \text{ or } \dots \text{ or } \{c_1, \dots, c_k\}$$

must also occur in that model⁸.

Given this, suppose that a model labelled M is detected by the shallow parser. Then we can check if this property (Exig_i) is verified, by counting the number of occurrences of a, b and c in the model labelled M. That is:

Let $\text{count}(x, X)$ be a function returning the number of occurrences of x in (the model) X. If

$$\begin{aligned} & \text{count}(a_1, M) \neq 0 \text{ and } \dots \text{ and } \text{count}(a_n, M) \neq 0 \\ & \text{and} \\ & [(\text{count}(b_1, M) = 0 \text{ or } \dots \text{ or } \text{count}(b_m, M) = 0) \\ & \text{and } \dots \text{ and} \\ & (\text{count}(c_1, M) = 0 \text{ or } \dots \text{ or } \text{count}(c_k, M) = 0)] \end{aligned}$$

then M does not satisfies Exig_i .

If the model M does not satisfies Exig_i , then it is marked with the identification of Exig_i .

As a predicate “count” is available in XSLT, mapping the previous formulas in XSLT is a trivial task.

Exigency properties – an example

As an example, consider an exigency property labelled Exig_{30} , defined as follows:

$$\text{Exig}_{30} \{a, b\} \Rightarrow_A \{c\}$$

It will be coded in XSLT as it follows:

⁸ These elements or others subsumed by them.

```

<xsl:template match='model[@name="A"]'>
<xsl:variable name="submodels" select="model"/>
<model>
<xsl:attribute name="name"><xsl:value-of select="@name"/>
</xsl:attribute>
<xml:if test="count($submodels[@name='a']!=0 and
count($submodels[@name='b']!=0 and
((count($submodels[@name='c']=0)))">
<exig>
<xsl:attribute name="id">E30</xsl:attribute>
</exig>
<xsl:if>

```

This XSLT file runs over a shallow parsed text and the segments not respecting this property are identified and marked. For example, if after a surface analysis we have the segment $(a\ b)_A$, that is:

```

<model name="A">
<model name="a"/>
<model name="b"/>
</model>

```

the transformation operated by Tira-Teimas returns:

```

<model name="A" exig="30">
<model name="a"/>
<model name="b"/>
</model>

```

which indicates that property number 30 was not respected.

In conclusion, we codify 5P properties in a stylesheet that will transform a XML text in another XML text where the sequences not verifying the 5P properties are marked.

Other properties

A similar approach is applied to exclusion properties. Uniqueness can also be treated through element's counting as in the following.

Consider the property

$$U_n = \{a_1, \dots, a_n\}.$$

If

$$\text{count}(a_1, M) \leq 1 \text{ and } \dots \text{ and } \text{count}(a_n, M) \leq 1$$

then the model M verifies this property. It should be noticed that a model with a coordination do not need to respect these properties.

4 Tira-Teimas Translator: an auxiliar tool

4.1 The problem

Writing 5P properties directly as a XSLT template is not an easy task. In fact, the complexity of the templates increases with the subsumption relations. The next example shows how difficult it is to write a basic 5P property as a XSLT template (in the following, `det` subsumes `artd_s`, `artd_p`, `dem` and `arti_s`).

The property

$$\text{Exig}_3 \{npr2\} \Rightarrow_{Nn} \{det\}$$

originates the following template:

```
<xsl:template match='model[name="m_nn"]'>
<xsl:variable name="submodels" select="model"/>
<model>
<xsl:attribute name="name">
  <xsl:value-of select="@name"/>
</xsl:attribute>

<xsl:if test="count($submodels[@name='npr2'])!=0 and (((
  count($submodels[@name='artd_s']) =0 and
  count($submodels[@name='artd_p']) =0 and
  count($submodels[@name='dem'])=0 and
  count($submodels[@name='arti_s'])=0 )))">
<exig>
  <xsl:attribute name="id">Exig3</xsl:attribute>
</exig>
</xsl:if>
```

4.2 Tira-Teimas Translator: generating Tira-Teimas

In order to simplify this situation, 5P properties are written in a friendlier XML format, and then an extra program – Tira-Teimas Translator (TTT) – maps these properties into XSLT. That is, TTT generates Tira-Teimas (see Fig. 2).

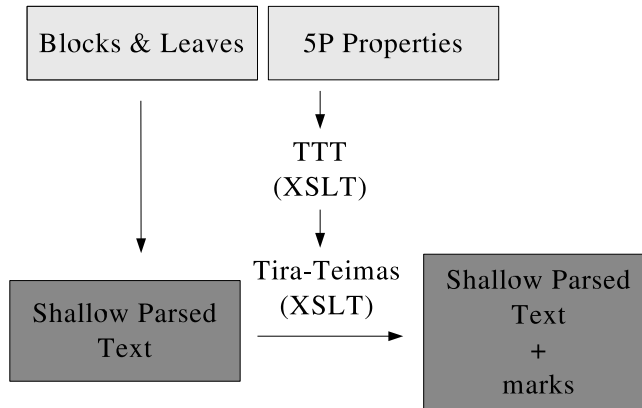


Figura 2. Tira-Teimas Translator

4.3 An example

In order to implement TTT we also used XSLT. In the following the property presented before is coded in the new XML format:

```

<exig m_label="m_nn" num="Exig3">
<left_exig><element id="npr2"/></left_exig>
<right_exig><element id="det"/></right_exig>
</exig>
  
```

TTT picks up the properties written in this format and returns Tira-Teimas.

5 Results

We took several (non treated) corpora, such as a tale, a thesis' chapter, a text about meteorology and others and we parsed them with Susana. Also, we coded the set of 5P properties describing nominal phrases, presented in [3], in XML, in the friendlier format preciously discussed. Then, TTT run over the 5P properties in XML, generating Tira-Teimas. Finally, Tira-Teimas checked the shallow parsed corpora. In the following we show the most interesting situations we got.

Consider the exigency property number 15⁹:

Exig₁₅: adj_s ⇒ det | algum | qualquer | ... | tanto

⁹ *algum*, *qualquer*, ..., *tanto* are very particular category labels for the words *algum*, *qualquer*, ..., *tanto*, respectively.

Consider also the following two hypothesis obtained by Susana:

$$\begin{aligned} & (muito_{q3-s} \ alto_{adj1-s})_{Nn} \\ & (((muito_{adv21_1})_{Advn} \ alto_{adj1-s})_{An1})_{Nn} \end{aligned}$$

Tira-Teimas detected that the first hypotheses was not respecting property Exig₁₅, as none of the elements needed by *adj1_s* (subsumed by *adj_s*) are in the sequence. So, we are now able to cut the first hypotheses, solving an ambiguity situation.

Consider now the following property:

$$\text{Exig}_5: \text{npr1} \Rightarrow_{Nn} \text{NIL}$$

and Susana's result:

$$\begin{aligned} & (a_{art-s} \ \text{Dublin}_{npr1})_{Nn} \\ & (a_{prep} \ \text{Dublin}_{npr1})_{Pn} \end{aligned}$$

According to Tira-Teimas, the first hypotheses goes against property Exig₅, as inside an *Nn*, an *npr1* should be alone. Once again, an ambiguity situation can be eliminated.

Now, take the exclusion property number 25:

$$\text{Exclu}_{25}: \text{ppas} \not\Rightarrow \text{adj3}$$

Susana returned the following sequences:

$$\begin{aligned} & (um_{artd-p} \ \text{comprimido}_{ppas} \ \text{vermelho}_{adj3-p})_{Nn} \\ & (um_{artd-p} \ \text{comprimido}_{nc1-s} \ \text{vermelho}_{adj3-p})_{Nn} \end{aligned}$$

As within an *Nn*, *adj3* excludes *ppas*, the first sequence is not possible.

To conclude, consider the property:

$$\text{Exig}_4: \text{npr4} \Rightarrow_{Nn} \text{det} \mid \text{algum} \mid \text{qualquer} \mid \dots \mid \text{card}_p$$

and the following result:

$$(\text{Alfredo}_{npr4})_{Nn} \ \text{teve...}$$

if this property is correct an *npr4* (proper nouns for people's names), can not appear alone within an *Nn*. As this is a possible situation, we can conclude that this property is not correct.

6 Conclusions and further work

Tira-Teimas guarantees that models not respecting 5P properties are identified. As so, bad results can be eliminated. Also, bad properties can be eliminated.

Although the original motivation for Tira-Teimas was to check 5P properties, Tira-Teimas is not bounded to this application. In fact, it can also be used to find differences between what is syntactically correct and what is currently practiced.

In addition, Tira-Teimas could be applied to detect differences between the Portuguese from Portugal and Portuguese from Brazil. For example, the 5P properties from [3] describing Portuguese (from Portugal) nominal phrases could be applied to a Brazilian shallow parsed text.

Another possible application to Tira-Teimas would be to retrieve the 5P properties from a disambiguated shallow parsed text. All the possible 5P properties would be automatically generated and all the properties that the shallow parsed text disrespects would be possible restrictions over the language syntax.

Referências

1. Bès, G.G.: La phrase verbal noyau en français. In: in Recherches sur le français parlé, 15, Université de Provence, France (1999) 273–358
2. Bès, G.G., Hagège, C.: Properties in 5P (soon in the GRIL pages). (November, 2001)
3. Hagège, C.: Analyse Syntactic Automatique du Portugais. PhD thesis, Université Blaise Pascal, Clermont-Ferrand, France (2000)
4. Bès, G.G., Hagège, C., Coheur, L.: Des propriétés linguistiques à l'analyse d'une langue. In: VEXTAL, Venice, Italy (1999)
5. Batista, F., Mamede, N.: SuSAna: Módulo multifuncional da análise sintáctica de superfície. In Gonzalo, J., Peñas, A., Ferrández, A., eds.: Proc. Multilingual Information Access and Natural Language Processing Workshop (IBERAMIA 2002), Sevilla, Spain (2002) 29–37
6. World Wide Web Consortium (W3C): (Extensible Markup Language (XML)) See: www.w3.org/XML.
7. World Wide Web Consortium (W3C): (The Extensible Stylesheet Language (XSL)) See: www.w3.org/Style/XSL.